

FINAL PROJECT

Presented by:

Kelompok 12

Andi Dhiya Rachmy Ariffia

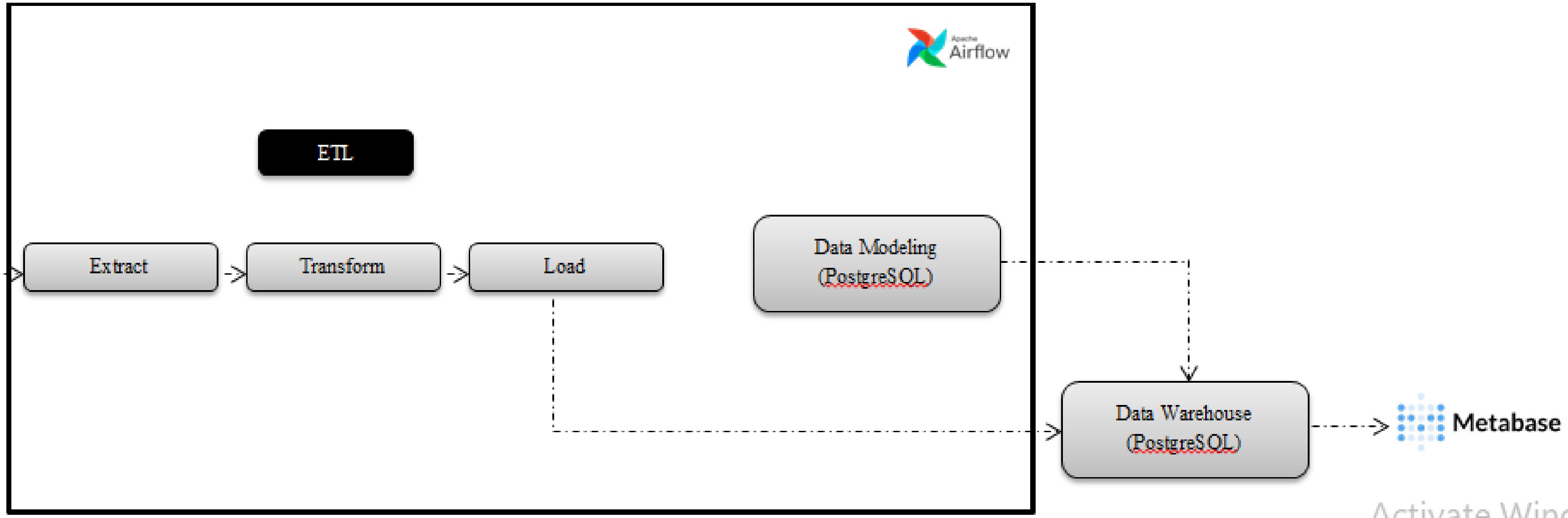
Haidar Fadhila Fiqa

Meiyra Istiqomah

Dita Artiningtyas

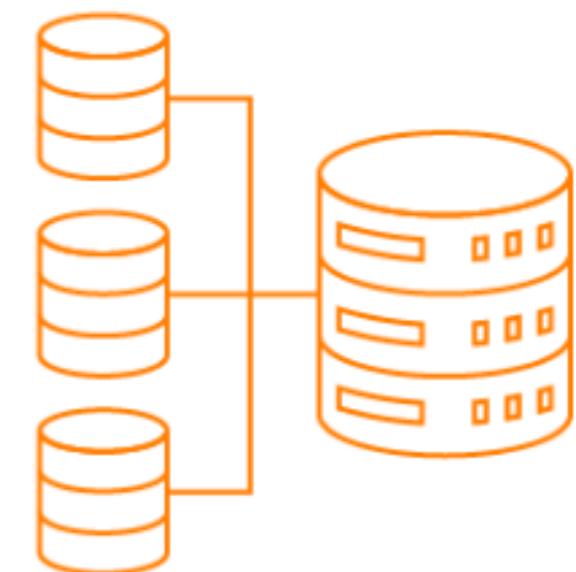
Mentor Pendamping: Ardhi Wahyudi

ARSITEKTUR



ETL

ETL stands for Extract, Transform, Load, and it refers to a process in data management and data warehousing.



Extract

Retrieves and verifies data from various sources



Transform

Processes and organizes extracted data so it is usable



Load

Moves transformed data to a data repository

PROSES ETL

Extract Function dan Transform

```
1 # Step 2: Extract Function
2 def extract():
3     path = "/opt/airflow/data"
4     onlyfiles = [f for f in listdir(path) if isfile(join(path, f))]
5
6     # now do the for loop
7     for i in onlyfiles:
8         #customer (csv)
9         if ".csv" in i and "customer" in i:
10             df_customer = pd.DataFrame()
11             for j in range(10):
12                 df = pd.read_csv(path+"/"+i)
13                 df_customer = pd.concat([df_customer,df])
14
15         #product dan product category (xls)
16         if ".xls" in i and "product" in i:
17             if "product_category" in i:
18                 df_product_category = pd.read_excel(path+"/"+i)
19             else:
20                 df_product = pd.read_excel(path+"/"+i)
21
```

The extract function reads data from different file formats (CSV, XLS, JSON, Parquet, Avro) located in the specified directory (/opt/airflow/data). The extracted data is then returned as a list of Pandas DataFrames.

PROSES ETL

Import Modul

```
1 # Step 1: Importing Modules
2 from airflow import DAG
3 from datetime import timedelta, datetime
4 from airflow.operators.python_operator import PythonOperator
5 import pandas as pd
6 import fastavro
7 import psycopg2
8 from os import listdir
9 from os.path import isfile, join
10 from sqlalchemy import create_engine
```

PROSES ETL

LOAD FUNCTION

```
1 # Step 3: Load Function
2 def load(*args, **kwargs):
3     connection = psycopg2.connect(
4         host="dataeng-warehouse-postgres",
5         port=5432,
6         dbname="data_warehouse",
7         user="user",
8         password="password"
9     )
10
11    cursor = connection.cursor()
12    list_tabel = extract()
13    tabel = {}
14    tabel['customers'] = list_tabel[0]
15    tabel['products'] = list_tabel[1]
16    tabel['product_categories'] = list_tabel[2]
17    tabel['suppliers'] = list_tabel[3]
18    tabel['login_attempt_history'] = list_tabel[4]
19    tabel['coupons'] = list_tabel[5]
20    tabel['orders'] = list_tabel[6]
21    tabel['order_items'] = list_tabel[7]
22
23    engine = create_engine('postgresql://user:password@dataeng-warehouse-postgres:5432/data_warehouse')
24
```

The load function establishes a connection to a PostgreSQL database (data_warehouse) and loads the extracted data into corresponding tables.

PROSES ETL

CREATING TASKS

```
1 # Step 5: Creating Tasks
2 task_1 = PythonOperator(
3     task_id='task_1',
4     python_callable=extract,
5     provide_context=True,
6     dag=dag
7 )
8
9 task_2 = PythonOperator(
10    task_id='task_2',
11    python_callable=load,
12    provide_context=True,
13    dag=dag
14 )
```

Two PythonOperator tasks (task_1 and task_2) are created. task_1 calls the extract function, and task_2 calls the load function.

SETTING UP DEPENDENCIES

```
1 # Step 6: Setting up Dependencies
2 task_1 >> task_2
```

This step defines the dependency between task_1 and task_2, indicating that task_2 should only run after task_1 has completed successfully.

DATA MODELLING

Data modelling digunakan untuk menyusun dan menyimpan data dalam bentuk yang memfasilitasi analisis bisnis yang cepat dan efisien

dim_product_supplier

Terdiri dari table :

1. Product
2. Product Category
3. Supplier

dim_order_coupon

Terdiri dari table :

1. Order
2. OrderItem
3. Coupons

dim_orders

Terdiri dari table Orders

dim_customers_login

Terdiri dari table :

1. Customers
2. Login attempt history



DATA MODELLING

Fact Table Login History

Fact table tersebut berisi informasi kapan customer melakukan log in, dikelompokkan sebagai berikut :

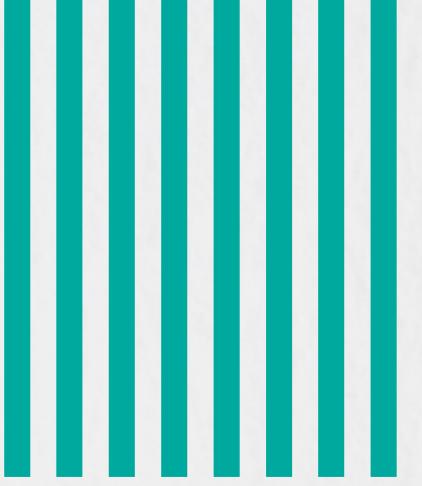
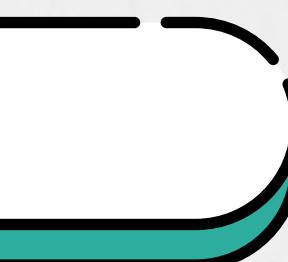
- a. Dini Hari
- b. Pagi
- c. Siang - Sore
- d. Malam

Fact Subtotal Diskon

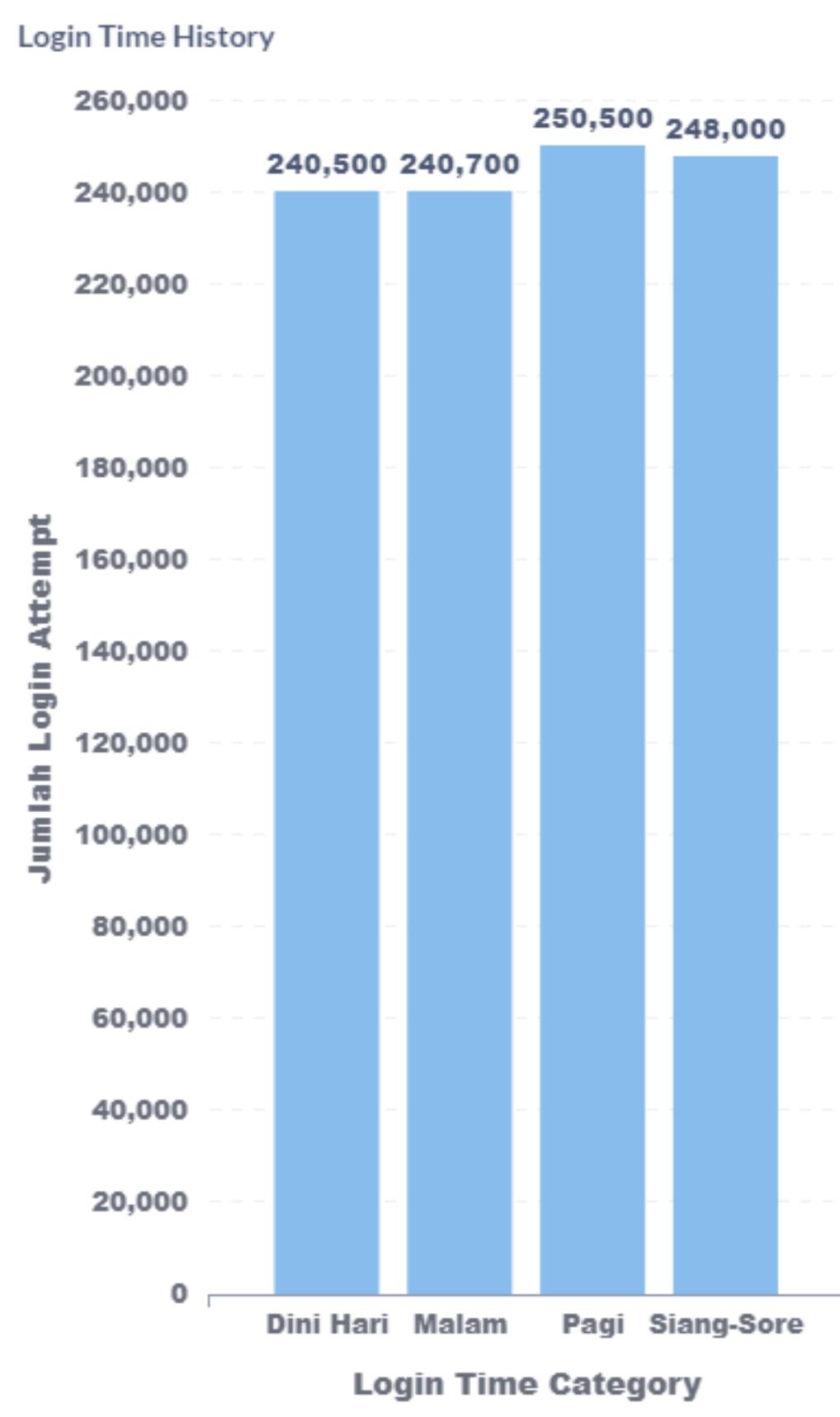
Fact table tersebut berisi informasi total biaya yang dikeluarkan oleh customer setelah memperhitungkan jumlah diskon yang didapatkan

Fact Table Supplier based on Product Category

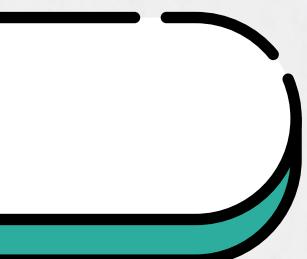
Fact table tersebut berisi informasi terkait jumlah supplier untuk setiap product category



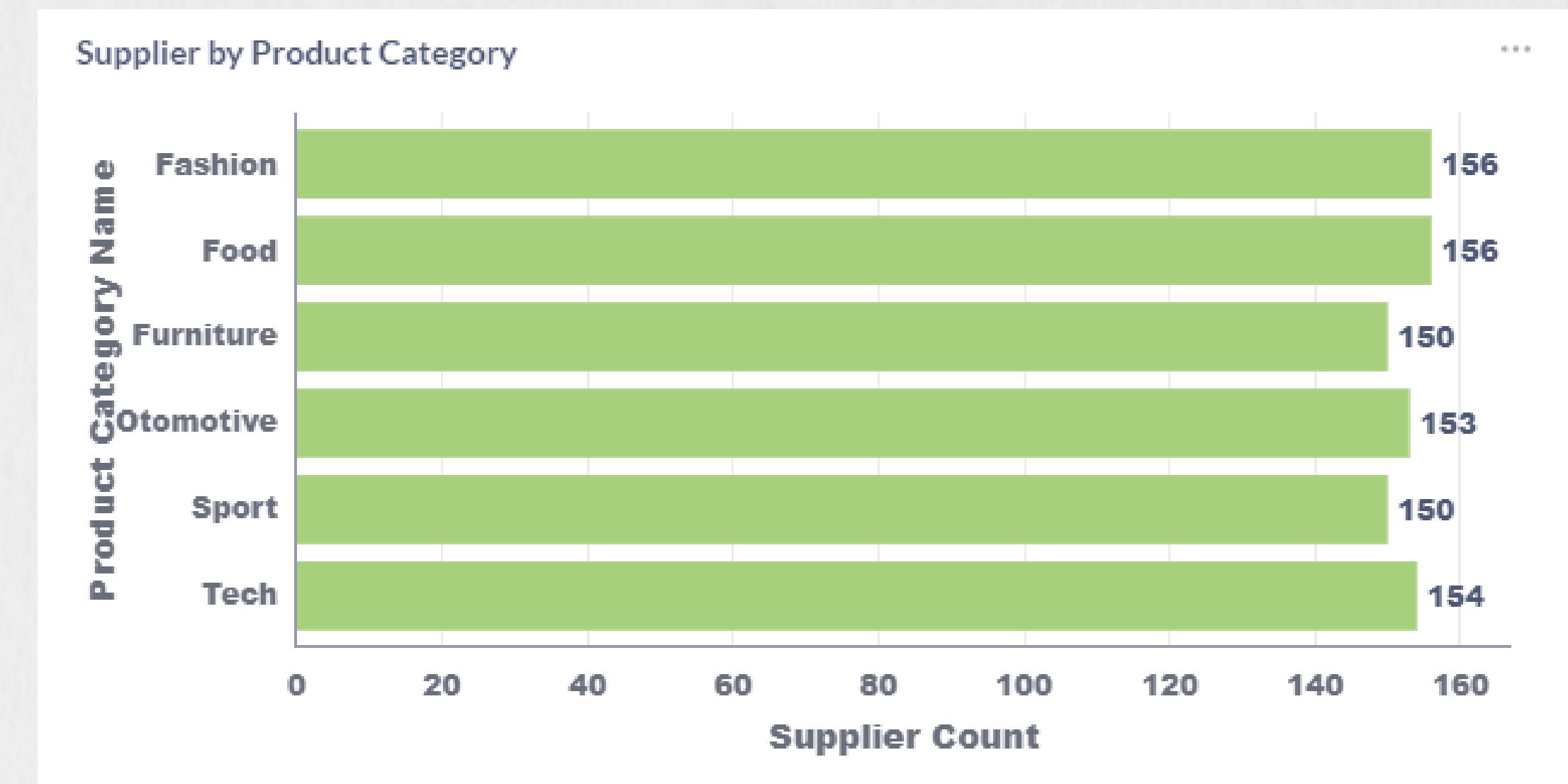
VISUALISASI



Berdasarkan hasil visualisasi, dapat diketahui bahwa customer banyak melakukan login pada pagi hari, yaitu sebanyak 250.500 customer.



VISUALISASI

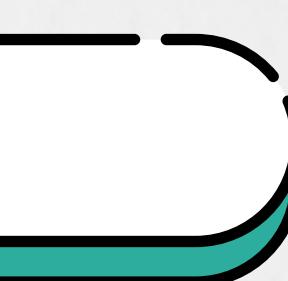


Berdasarkan hasil visualisasi, dapat diketahui bahwa kategori “Fashion” dan “Food” merupakan product category yang memiliki supplier paling banyak, yaitu sebanyak 156 suppliers.

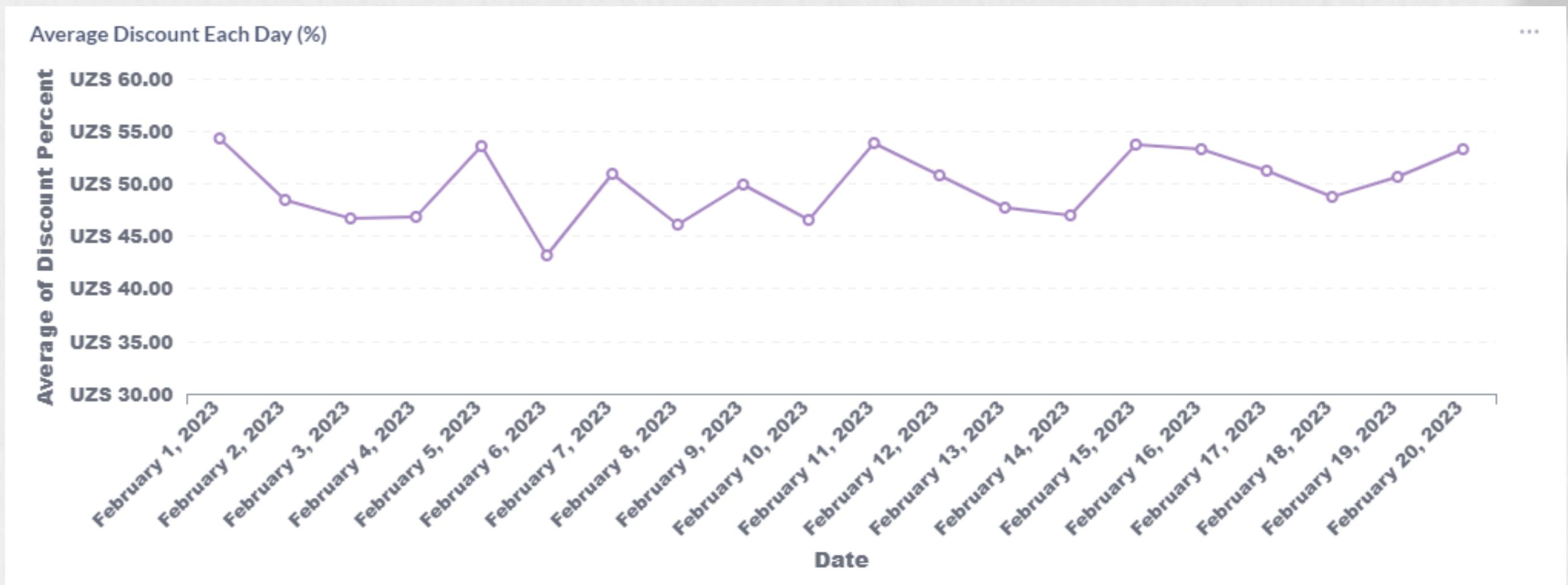
VISUALISASI

Country	Count
Albania	1
Algeria	1
American Samoa	1
Antarctica (the territory South of 60 deg S)	2
Antigua and Barbuda	1
Australia	1
Austria	1
Azerbaijan	1
Bahrain	1

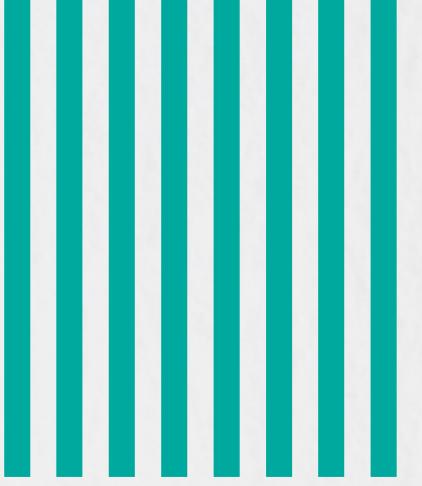
Visualisasi diatas merupakan visualisasi dalam bentuk table yang berisi tentang persebaran wilayah para supplier, yang mana sebagian besar terdiri 1 supplier setiap daerahnya. Namun, ada beberapa daerah yang juga memiliki banyak supplier, misalnya seperti Antarctica yang memiliki 2 supplier dan Timor Leste memiliki 3 supplier.



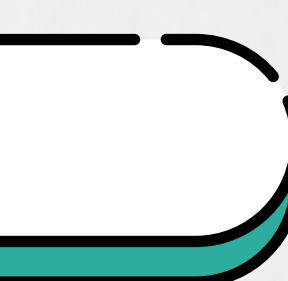
VISUALISASI



Berdasarkan hasil visualisasi, dapat diketahui rata rata diskon setiap harinya. Terlihat bahwa pada 11 Februari 2023 memiliki rata-rata diskon yang paling besar, yaitu sebesar 54%



SARAN BISNIS

- 
1. Optimalkan pelayanan pada pagi hari. Manfaatkan peluang dengan menyediakan pelayanan atau penawaran spesial pada pagi hari untuk meningkatkan keterlibatan pelanggan.
 2. Pertimbangkan strategi pemasaran untuk kategori produk unggulan pada kategori "Fashion" dan "Food" yang memiliki jumlah supplier paling banyak. Evaluasi dan tingkatkan kerjasama dengan supplier di kategori ini.
 3. Analisis lebih lanjut tentang wilayah dengan banyak supplier. Fokuskan analisis lebih lanjut pada wilayah-wilayah yang memiliki lebih dari satu supplier untuk mengidentifikasi peluang bisnis atau potensi peningkatan efisiensi operasional



KESIMPULAN

Project yang dilakukan telah berhasil menggabungkan proses ETL, pemodelan data dan analisis visualisasi dengan efektif. Pada pemodelan data memberikan dasar yang kuat untuk analisis, sementara visualisasi dengan metabase bisa memberikan penjelasan yang jelas





THANK
YOU

